# Data Analytics at NERSC

**Joaquin Correa**
**JoaquinCorrea@lbl.gov**
**NERSC Data and Analytics Services**

**NERSC User Meeting**
**August, 2015**

# Data analytics at NERSC

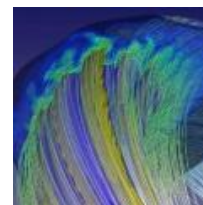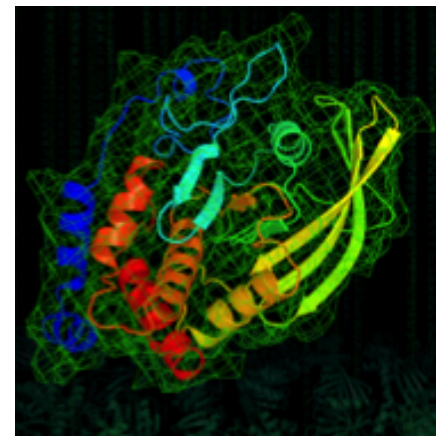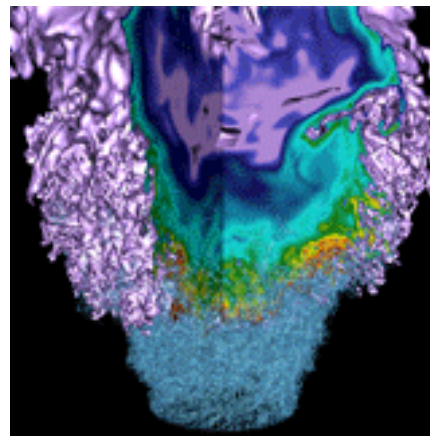| | | | | |
|---|---|---|---|---|
| **Science Applications** | Climate, Cosmology, Kbase, Materials, BioImaging, Your science! | | | |
| **Analytics Capabilities** | Statistics, Machine Learning | Image Processing | Graph Analytics | Database Operations |
| **Tools + Libraries** | R, python, MLBase | MATLAB OMERO, Fiji | GraphX | SQL |
| **Runtime Framework** | MPI | Spark | SciDB | |
| **Resource Management** | Filesystems (Lustre), Batch/Queue Systems | | | |
| **Hardware** | SandyBridge/KNL chipset, Burst Buffers, Aries Interconnect | | | |

# Data analytics at NERSC

| Analytics Capabilities | Statistics, Machine Learning | Image Processing | Graph Analytics | Database Operations |
| --- | --- | --- | --- | --- |
| Tools + Libraries | R, python, MLBase | MATLAB OMERO, Fiji | GraphX | SQL |

| Runtime Framework | MPI | Spark | SciDB |
| --- | --- | --- | --- |

# Talk Overview

- **Data analytics tools**
- **Data insight**
- **Scale your analysis**

# Talk Overview

- **Data analytics tools**
- **Data insight**
- **Scale your analysis**

- R is a language and environment for statistical computing and graphics. It provides a wide variety of statistical tools, such as linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, graphics, and it is highly extensible.

```r
# Goal: To do `moving window volatility' of returns.

library(zoo)

# Some data to play with (Nifty on all fridays for calendar 2004) --
p <- structure(c(1946.05, 1971.9, 1900.65, 1847.55, 1809.75, 1833.65, 1913.6, 1852.65

# Shift to returns --
r <- 100*diff(log(p))
head(r)
summary(r)
sd(r)

# Compute the moving window vol --
vol <- sqrt(250) * rollapply(r, 20, sd, align = "right")

# A pretty plot --
plot(vol, type="l", ylim=c(0,max(vol,na.rm=TRUE)),
     lwd=2, col="purple", xlab="2004",
     ylab=paste("Annualised sigma, 20-week window"))
grid()
legend(x="bottomleft", col=c("purple", "darkgreen"),
       lwd=c(2,2), bty="n", cex=0.8,
       legend=c("Annualised 20-week vol (left scale)", "Nifty (right scale)"))
par(new=TRUE)
plot(p, type="l", lwd=2, col="darkgreen",
     xaxt="n", yaxt="n", xlab="", ylab="")
axis(4)
```
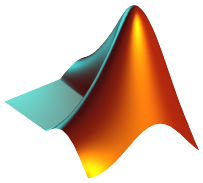
- MATLAB is a technical computing language that integrates computation, visualization, and programming in an easy-to-use environment where problems and solutions are expressed in familiar mathematical notation.

Toolboxes:

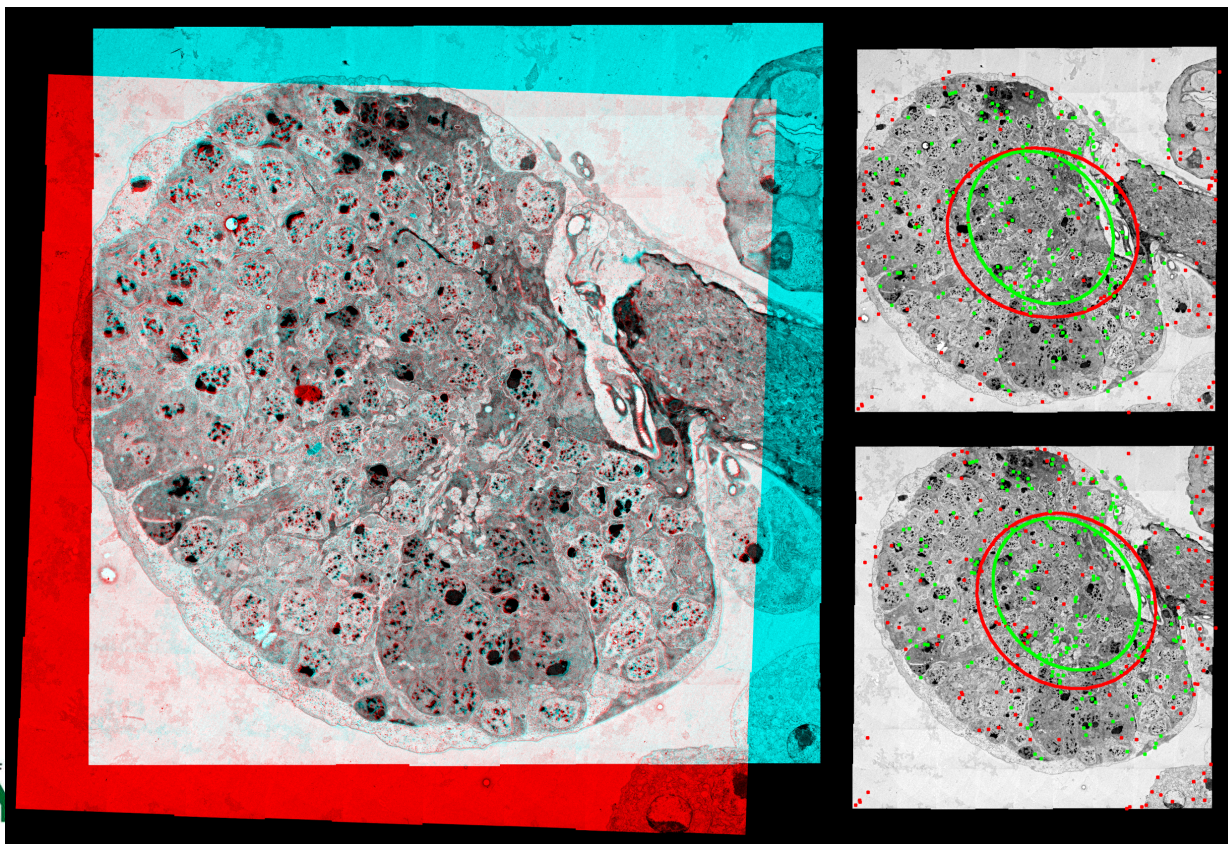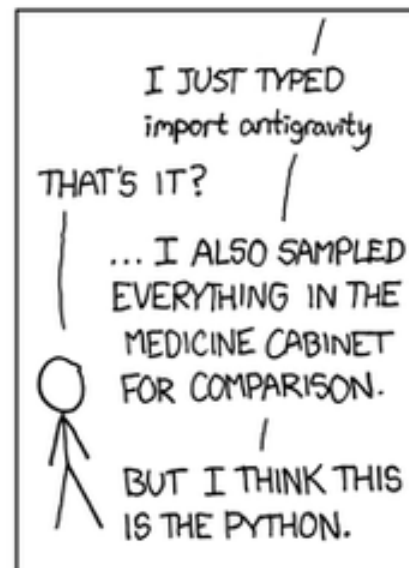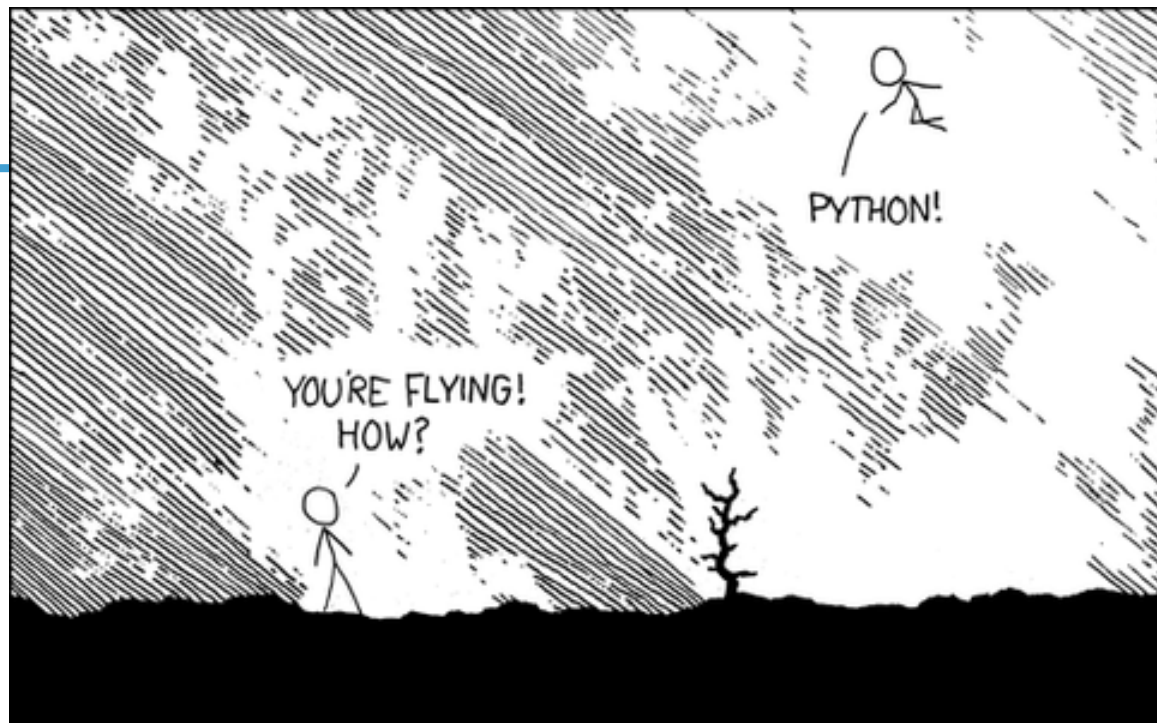| MATLAB | 16 |
|---|---|
| Image Processing | 2 |
| Neural networks | 1 |
| Optimization | 2 |
| Parallel computing | 2 |
| Signal processing | 1 |
| Statistics | 2 |
| Compiler | 1 |

# Mathematica

- Mathematica is a fully integrated environment for technical computing. It performs symbolic manipulation of equations, integrals, differential equations, and most other mathematical expressions. Numeric results can be evaluated as well.

- Fiji Is Just ImageJ - Fiji is an image processing package. It can be described as a "batteries-included" distribution of ImageJ (and ImageJ2), bundling Java, Java3D and a lot of plugins organized into a coherent structure.
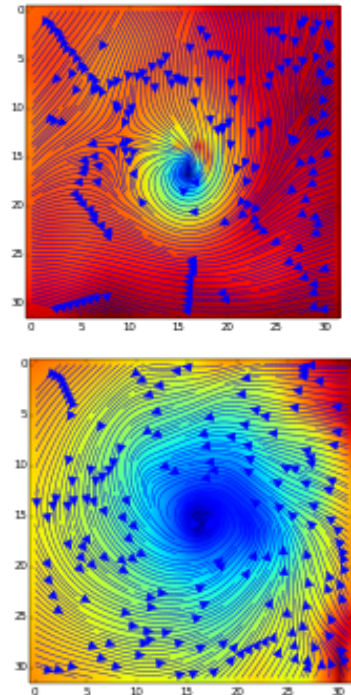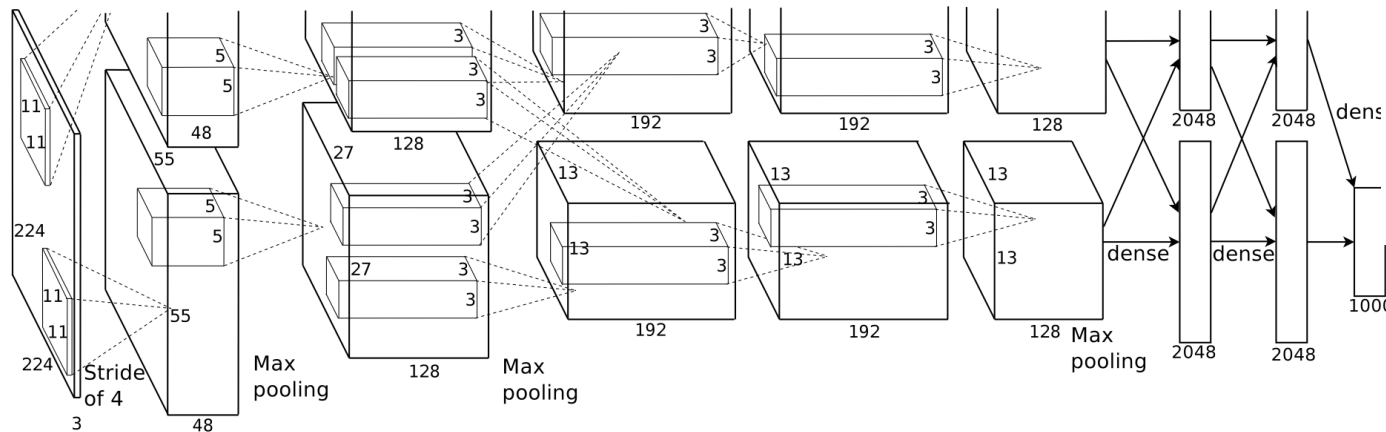
- Scientific Computing Tools for Python
  - NumPy
  - The SciPy library
  - Matplotlib
  - pandas
  - SymPy
  - IPython
  - nose
  - Cython
  - Scikits
  - h5py
  - mpi4py

# Deep learning at NERSC

- <u>neon</u> is an easy to use, python-based scalable Deep Learning library.

Deep Learning has recently achieved state-of-the-art performance in a wide range of domains including images, speech, and text. It is seeing adoption in the HPC community as a tool for large-scale data processing.
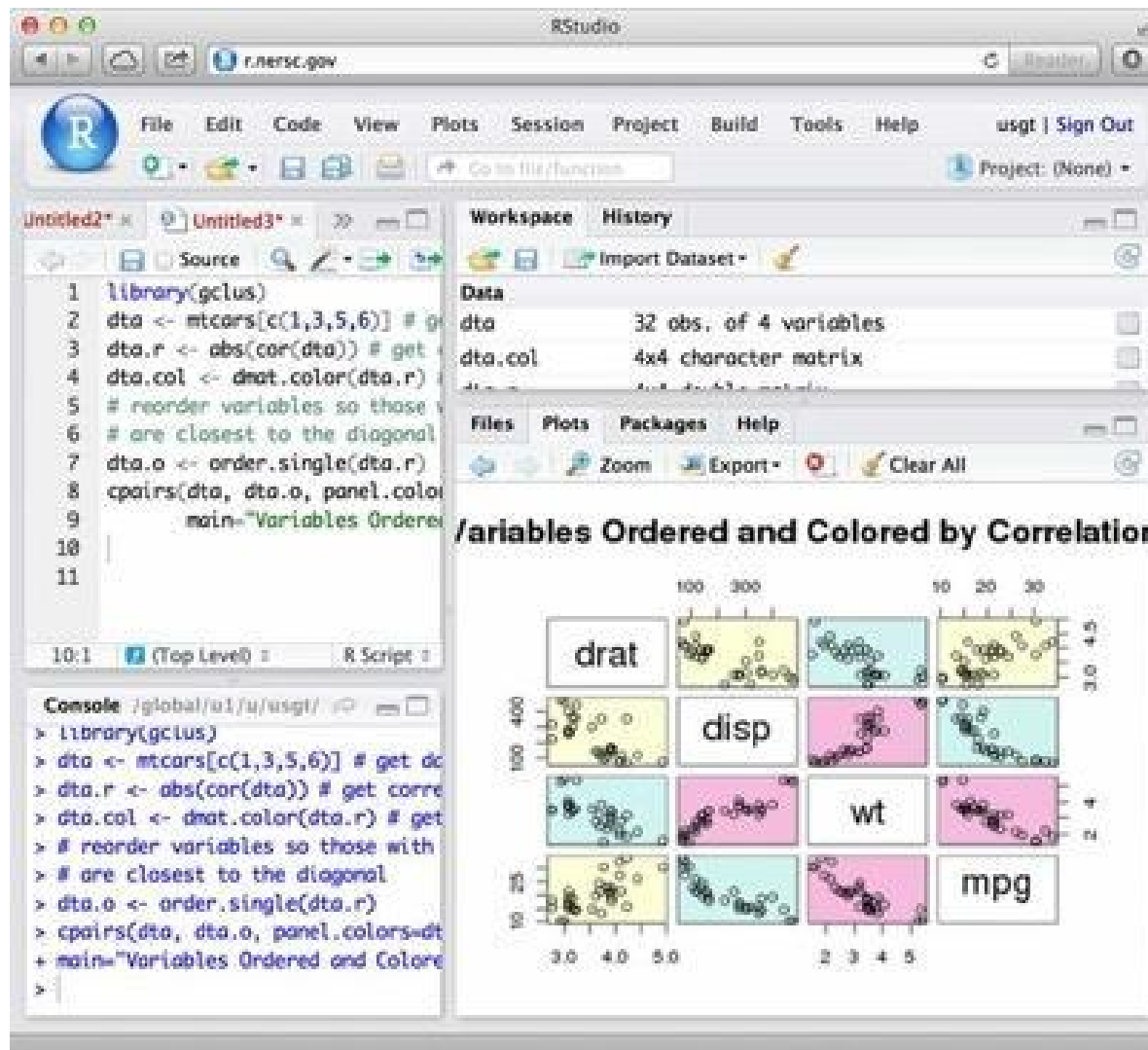
# Talk Overview

- Data analytics tools
- **Data insight**
- Scale your analysis

# R and RStudio Service(Beta)

# iPythonNotebook Service(Beta)

# Talk Overview

- Data analytics tools
- Data insight
- **Scale your analysis**

- Apache Spark is a fast and general-purpose cluster computing system. It provides high-level APIs in Java, Scala, Python and R, and an optimized engine that supports general execution graphs.

- Easy to Prototype
  - Interactive shells make it easy to explore your data
  - Interactively debug your analyses
  - Works with iPython notebooks

- Easy to Run
  - High-Level API using map-reduce paradigm
  - Implement all your analyses in a few lines of
  - Python
  - Scala
  - Java

| Computation Type | Spark Implementation |
|---|---|
| Machine Learning | MLlib, Spark ML |
| Graph Computations | GraphX |
| Database Operations | Spark SQL |
| Streaming Analysis | Spark Streaming |
| Your Own Custom Analysis | Using Spark's Built In Functions |



Spark SQL | Spark Streaming | MLlib (machine learning) | GraphX (graph)

Apache Spark

Computation types can be combined seamlessly all in the same piece of code!

# SciDB For High Usability Big Data Analytic

- **Why?** It's painful to manage and analyze terabytes of data. Need a unified solution that's easy to use.

- **What?** SciDB is a parallel database for array-structured data, great for **Terabytes** of:
  - **Time series, spectrums, imaging**, etc

- The greatest benefit of SciDB is:
  - **Usability**: Use HPC hardware without learning parallel programming and parallel I/O.



SciDB Distribute a big array on many nodes

# NERSC Data Analytic Services



**Big and Diverse Computing Facility**

6000+ Users, 700+ Projects

3+ PetaFlops (20+pf more coming)

50+ PB Storage
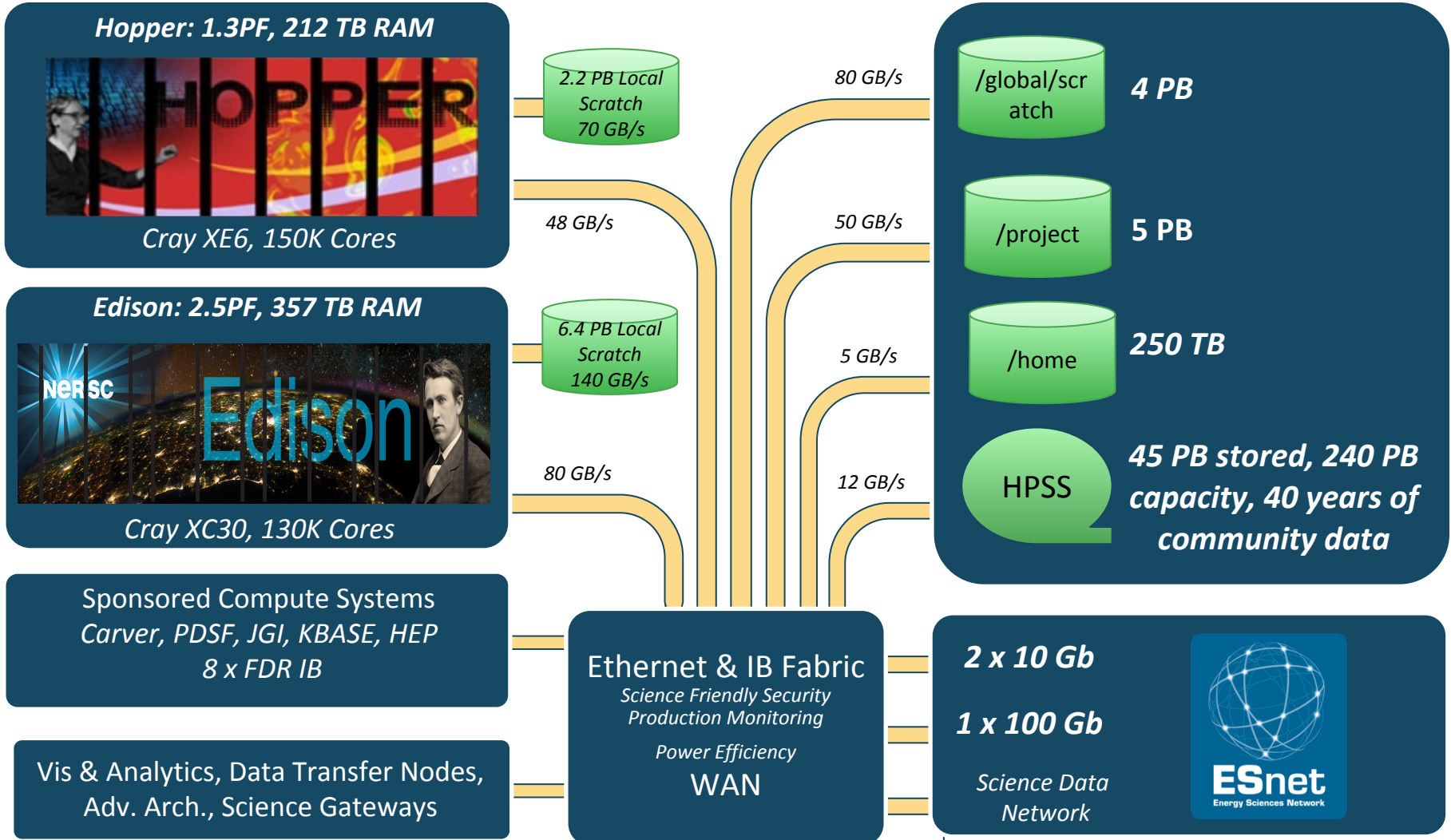


Production Data Services

Science Engagement

**Thank you.**

# Cori: Unified architecture for HPC and Big Data

- **64 Cabinets of Cray XC System**
  - 50 cabinets 'Knights Landing' *manycore* compute nodes
  - 10 cabinets 'Haswell' compute nodes for *data partition*
  - ~4 cabinets of Burst Buffer
  - 14 external login nodes
  - Aries Interconnect (same as on Edison)
- **Lustre File system**
  - 28 PB capacity, 432 GB/sec peak performance
- **NVRAM "Burst Buffer" for I/O acceleration**
- **Significant Intel and Cray application transition support**
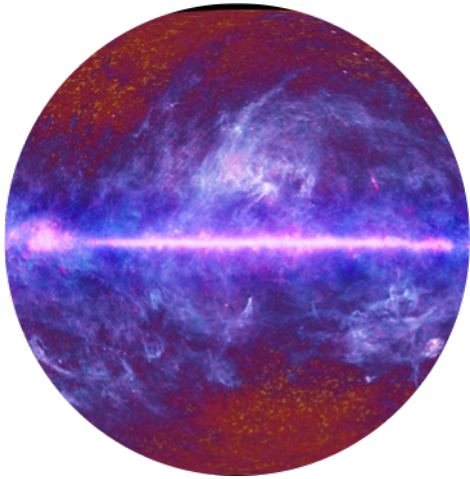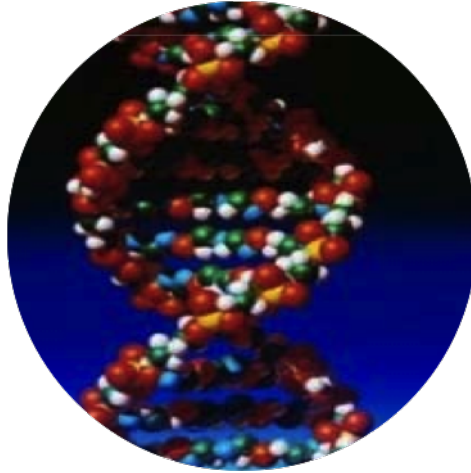- **Delivery in mid-2016; installation in new LBNL CRT**

# NERSC Systems



**Hopper: 1.3PF, 212 TB RAM**

*Cray XE6, 150K Cores*

**Edison: 2.5PF, 357 TB RAM**

*Cray XC30, 130K Cores*

Sponsored Compute Systems
*Carver, PDSF, JGI, KBASE, HEP*
*8 x FDR IB*

Vis & Analytics, Data Transfer Nodes,
Adv. Arch., Science Gateways

2.2 PB Local
Scratch
70 GB/s

6.4 PB Local
Scratch
140 GB/s

48 GB/s

80 GB/s

/global/scratch — *4 PB*

/project — *5 PB*

50 GB/s

5 GB/s

/home — *250 TB*

12 GB/s

HPSS — *45 PB stored, 240 PB capacity, 40 years of community data*

Ethernet & IB Fabric
*Science Friendly Security*
*Production Monitoring*

*Power Efficiency*
WAN

*2 x 10 Gb*

*1 x 100 Gb*

*Science Data Network*

ESnet
Energy Sciences Network

# 5 V's of Scientific Big Data

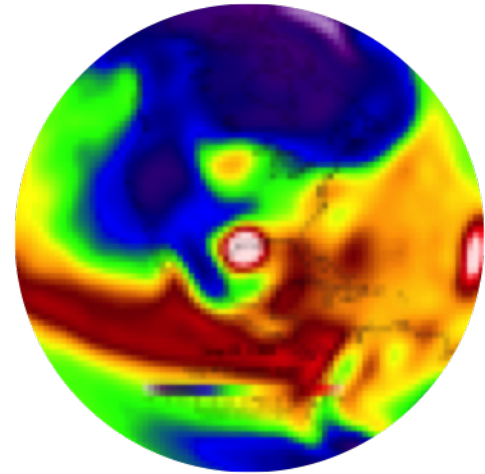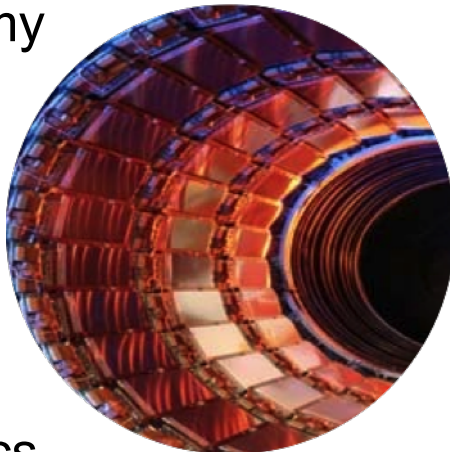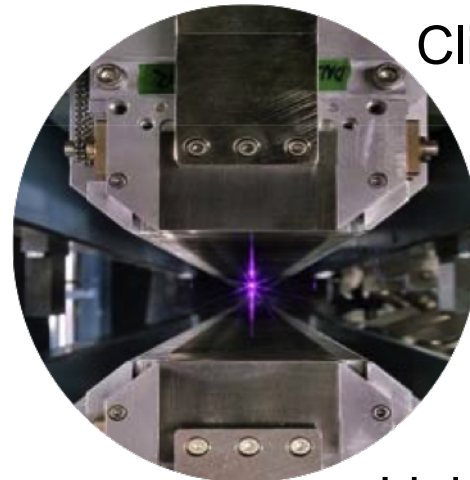| Science Domain | Variety | Volume | Velocity | Veracity |
|---|---|---|---|---|
| Astronomy | Multiple Telescopes, multi-band/spectra | O(100) TB | 100 GB/night – 10 TB/night | Noisy, acquisition artefacts |
| Light Sources | Multiple imaging modalities | O(100) GB | 1 Gb/s-1 Tb/s | Noisy, sample preparation/acquisition artefacts |
| Genomics | Sequencers, Mass-spec, proteomics | O(1-10) TB | TB/day | Missing data, errors |
| HEP: LHC, Daya Bay | Multiple detectors | O(100) TB – O(10) PB | 1-10 PB/s reduced to GB/s | Noisy, artefacts, spatio-temporal |
| Climate | Simulations Multi-variate, spatio-temporal | O(10) TB | 30-70 GB/s | 'Clean', need to account for multiple sources of uncertainty |

# DOE Facilities are Facing a Data Deluge



Astronomy

Genomics

Climate

Physics

Light Sources

- 26 -